

**CONTENT VALIDITY AND ITEM ANALISYS
OF SEMESTER II ENGLISH FINAL TEST FOR TENTH
GRADE STUDENTS OF SMAN3 SIDOARJO**

THESIS

**Submitted in Partial Fulfillment of the Requirements
for the Degree of Sarjana Pendidikan Islam (S.Pd.I)
in Teaching English**

PERPUSTAKAAN IAIN SUNAN AMPEL SURABAYA	
No. KLAS T-2010 036 PRI	No. REG : T-2010 / PRI / 036 ASAL BUKU : TANGGAL :

Oleh :

MILLATUL ISLAMIYAH
D05206055

ENGLISH EDUCATION DEPARTMENT
FACULTY OF TARBIYAH
STATE INSTITUTE FOR ISLAMIC STUDIES SUNAN AMPEL
SURABAYA
2010

APPROVAL SHEET

This thesis by :

Name : Millatul Islamiyah

NIM : D05206055

Title : **CONTENT VALIDITY AND ITEM ANALYSIS ON SEMESTER II
ENGLISH FINAL TEST FOR TENTH GRADE STUDENTS OF SMAN 3
SIDOARJO**

Has been approved by the advisor and could be proposed to fulfill the requirement for the
Graduate Degree of Sarjana Pendidikan in English Department of Tarbiyah Faculty

Surabaya, August 3rd, 2010

Advisor,



Dra. Irma Soraya, M. Pd
NIP. 196709301993032004

EXAMINERS APPROVAL SHEET

Thesis entitled:

Content Validity and Item Analysis of Semester II English Final Test for Tenth Grade Students of SMAN 3 Sidoarjo has been accepted and approved by the broad of examiners of English Department of Tarbiyah Faculty State Institute for Islamic Studies Sunan Ampel Surabaya.

Surabaya, August 20th, 2010



Dean of Tarbiyah Faculty,

Dr. H. Nur Hamim, M. Ag
NIP. 196203121991031002

The Board of Examiners
Advisor/Chair Person,

Dra. Irma Soraya, M. Pd
NIP. 196709301993032004

Secretary,

Siti Asmiyah, S. Pd
NIP. 197704142006642003

Examiner I,

Dr. Mohammad Salik, M. Ag
NIP. 196712121934031002

Examiner II,

Masdar Hilmy, M. A, Ph. D
NIP. 197103021996031002

SURAT PERNYATAAN

Saya yang bertanda tangan di bawah ini adalah:

Nama : Millatul Islamiyah
Nim : D05206055
Semester : VIII (Delapan)
Jurusan : Pendidikan Bahasa Inggris (PBI)
Fakultas : Tarbiyah
Alamat : Prastian Labuhan Sregeh Sampang Madura

Dengan ini menyatakan dengan sebenarnya bahwa skripsi yang berjudul "**Content Validity and Item Analysis of Semester II English Final Test for Tenth Grade Students of SMAN 3 Sidoarjo**", adalah asli dan bukan plagiat, baik sebagian maupun keseluruhannya.

Demikian pernyataan ini sesuai dengan sebenarnya, apabila pernyataan ini tidak sesuai dengan fakta yang ada, maka saya siap dimintai pertanggung jawaban sebagaimana peraturan perundang-undangan yang berlaku.

Surabaya, 3 Agustus 2010
Pembuat pernyataan

Millatul Islamiyah
NIM D05206055

ABSTRACT

Content Validity and Item Analysis of Semester II English Final Test for Tenth Grade Students of SMAN 3 Sidoarjo

Name : Millatul Islamiyah

NIM : D05206055

Advisor : Dra. Irma Soraya, M. Pd

Key Words: Content Validity, Item analysis, Index of Difficulty, Index of Discrimination, the Effectiveness of Distractors.

Testing is one kind of evaluation. As evaluation, testing is very needed to be applied in teaching to know the progress of the students. Without testing, the result of education will be foolish. In order to perform efficient and correctly, testing must be good designed. Nowadays, there are many teachers do not carefully writing a test. They ignore of the criteria of good test, they are content validity, reliability, index of difficulty, index of discrimination, and the effectiveness of distractors. This study is aimed to know what is Semester II English Final Test for Tenth Grade of SMAN 3 Sidoarjo like in term of content validity, index of difficulty, index of discrimination, and the effectiveness of distractors .

The design used in this study is a descriptive research because it describes the quality of Semester II English Final Test for Tenth Grade of SMAN 3 Sidoarjo. It also used quantitative approach since it used numerical calculation to compute the data. The object of this study is Semester II English Final Test for Tenth Grade of SMAN 3 Sidoarjo and only focus on multiple choice items, while the samples are X1, X2, and X3 class which are taken by random sampling.

The result of this study reported that Semester II English Final Test for Tenth Grade of SMAN 3 Sidoarjo has good content validity. It also reported that the index of difficulty of Semester II English Final Test for Tenth Grade of SMAN 3 Sidoarjo used for X1 and X2 are acceptable, but they are recognized easy test for X3. Besides, the index of discrimination of this test is for X1, satisfactory for X2, and malfunction for X3. Moreover, the Semester II English Final Test for Tenth Grade of SMAN 3 Sidoarjo has good distractors for X1, X2, and X3.

unlimited access to knowledge and information. In conclusion, having excellent English language ability can raise better chances of getting jobs that pays more and boarder your knowledge and networking.

Considering the importance of English as well as it's a lot of significant, our government through The Department National Education has determined English as one of compulsory content of curriculum taught at junior and senior high school.

“Kurikulum pendidikan dasar dan menengah wajib memuat: a. Pendidikan agama, b. pendidikan kewarganegaraan, c. bahasa, d. matematika, e. ilmu pengetahuan alam, f. ilmu pengetahuan sosial, g. seni dan budaya, h. pendidikan jasmani dan olahraga, i. keterampilan/kejuruan, j. muatan local”.²

Then, in regulation attachment of ministry of national education explained that language curriculum contains English language.³

English learning aimed in junior high school is oriented to reach functional level. It means that the students should be able to communicate oral and written in their daily life activity. While, English learning in senior high school is expected to reach informational level, because they have been prepared to continue their study in university.⁴ Nowadays English language is taught in elementary school even it is introduced earlier to children in the kindergarten. It is aiming to teach English language earlier in order to facilitate children to learn

² Peraturan Republik Indonesia Tentang Sistem Pendidikan Nasional No. 20 tahun 2003. Chapter X, section 37, verse 1.

³ Lampiran Peraturan Menteri Pendidikan Tentang Sistem Pendidikan Nasional No. 20 tahun 2003. Page 8.

Depdiknas, *Standard Kompetensi Mata Pelajaran Bahasa Inggris SMP dan Madrasah Tsanawiyah*, (Depdiknas: Jakarta, 2004) page 3

easier in the higher level of education. This phenomena show us several efforts conducted by government to provide generation with skillful.

In addition, to measure the success in teaching learning English, evaluation is needed evaluation is very important means in teaching learning process. It is used to know whether the teaching learning activities is success or not. The result conducted from evaluation provides wide information to the teachers to manage their classroom activities, as Dickins and Germaine say:⁵

“....it can provide a wealth of information to use for the future direction of classroom practice, for the planning of course, and for the management of learning task and students,”

Furthermore, evaluation and teaching is cannot be separated,⁶ all teaching process should be followed by evaluation indeed. Without evaluation, it seems impossible to measure as well as report students progress objectively.

Evaluation as Arikunto has said is the process of evaluating teaching learning process.⁷ There are several types of evaluation. One of them is test. Test is a series of question of measuring skill, knowledge, intelligences, and capacities of individual or group.⁸ Nevertheless it is common that sometimes evaluation considered has same meaning as testing, and that while students are being tested evaluation is taking place. However testing is only one component in the evaluation process.

⁵ Pauline Rea-Dickins& Kevin Germaine, *Evaluation*,(New York, Oxford University Press, 2008), Page 3

⁶ J. B Heaton, *Writing English Language Test*, (New York: Longman Group, 1988). Page 5

⁷ Suharsimi Arikunto, *Dasar-dasar Evaluasi Pendidikan*, (Jakarta, Bumi Aksara, 1993), Page

⁸ Ibid, Page 29

Furthermore, Muchtar Bukhori says:

“Tes ialah suatu percobaan yang diadakan untuk mengetahui ada atau tidaknya hasil-hasil pelajaran tertentu pada seseorang murid atau kelompok murid”³

In the other word, Kubizyn and Borich stated in their book, that test is just as tools that can contribute importantly to the process of evaluating pupils, the curriculum, and the teaching method.⁴

Those above are several definitions about test created by some experts. Although these were written in different words or sentences, however it expressed the same meaning that test is one tool of process evaluating pupils, curriculum, and teaching method to measure the skill, the work of curriculum, and the successful of the teaching method. In addition to the previous explanation, as one type of measurement, a test is necessarily quantifies characteristics of individuals to explicit procedures.⁵

³ Ibid. Page 35

⁴ Tom Kubiszyn and Gary Borich, *Educational Testing and Measurement* (Singapore, John Wiley & Sons, INC, 2003), Page 1

⁵ Lyle F. Bachman. *Fundamental Consideration in Language Testing*. (New York, Oxford University Press, 1990), Page 20

B. Purpose of Test

As stated in the previous chapter that test has interrelated with teaching as well as education. Language tests also have many uses in educational programs, and sometimes two or more purposes cover the same test.⁶

David conducted six objectives of language testing:⁷

1. To determine readiness for instructional programs.
2. To classify or place individuals in appropriate language classes.
3. To diagnose the individual's specific strengths and weaknesses.
4. To measure aptitude for learning.
5. To measure the extent of student achievement of the instructional goals.
6. To evaluate the effectiveness of instruction.

In addition, Arikuto said that testing has several purposes for education, such as: a) testing is able to select the good student, b) testing is able to diagnosing the strength and weakness of student, c) testing is able to place students in proper class that fits their ability, and d) testing is able to measure the effectiveness of the program employed.⁸

⁶ David P Harris, *Testing English as a second Language*, (New York: Mc Craw-ill, Inc, 1959), Page 2

⁷ Ibid

⁸ Suharsimi Arikunto, *Dasar-dasar Evaluasi Pendidikan*, (Jakarta: Bumi Aksara, 1993), Page

In similar words, several purposes of test as cited from Hughes as follows:⁹

1. To measure language proficiency
2. To discover how successful students have been in achieving the objective of a course of study
3. To diagnose students' strengths and weaknesses, to identify what they know and what they do not know.
4. To assist placement of students by identifying the stage or part of a teaching programme most appropriate to their ability.

Ebel Also mentioned several benefits of test for both students and teacher/instructor,¹⁰ such as, to measure students achievement and thus to contribute to the evaluation of educational progress and attainment. Test also benefits to motivate and direct students learning. In the generally fact, students tend to learn harder when they are examined or tested. They also stress to learn on the subjects that are tested. Not last, test can cause instructor/teacher to think carefully about the goals of instruction in a course.

⁹ Arthur Hughes, *Testing for Language Teachers*, (Australia, Cambridge University Press, 1989), Page 8

¹⁰ Robert L. Ebel, *Essential of Educational Measurement*. (USA, Prentice-Hall INC New Jersey, 1979). Page 22

In addition, Achievement test are divide into three types of test, they are:

1) Entry Behavior Test

Entry behavior test is held before student start to learn in the particular department of education (school). This kind of test is aimed to know whether the test taker provides abilities or skills which are required by particular school as accepting condition. It means those students who fulfill that criteria of accepting condition are able to accept in that school and vice versa.

2) Placement Test

Placement tests are used to give information that will help to put students at the stage of the teaching program most suitable to their abilities.¹⁷

For instance, an English course has three level of classroom, elementary, intermediate and advance. In order to put the students in the proper class based on their ability, diagnostic test is held. The criteria of score in each class have been determined. Students who get particular determined score for elementary level are proper to place in

¹⁷ Arthur Hughes, *Testing for Language Teacher*, ((Australia, Cambridge University Press, 1989), Page 17

b. Formative Test

Formative test is held during the activities of teaching learning is going on. Usually it is done in the end of the accomplishment of one course. Thus, formative test might be held many times in one semester.

Formative test is aimed when teachers need to check on the progress of their students, to see how far they have achieved what they should have learned.¹⁹ The information conducted from this kind of test enable teachers to measure the effectiveness of their classroom activities, and also to modify their future teaching plans.

c. Summative Test

Summative test is used at, say, the end of the term, semester or year. It is aimed to measure what has been achieved both group and individuals.²⁰ The materials that are tested including all course objectives which have been learned during one semester or year.

Unlike Formative test, summative test has more general purpose. The general purpose here have already stated clearly in the *standar isi* of every lesson.

¹⁹ Arthur Hughes, *Testing for Language Teachers*, (Australia, Cambridge University Press, 1989), Page 5

1989), Page 5

e. Aptitude tests.

An aptitude test serves to indicate an individual's facility for acquiring specific skills and learning.²⁴

Language learning aptitude is a complex matter, consisting of such factors as intelligence, age, motivation, memory, phonological sensitive and sensitivity to grammatical patterning.²⁵

Aptitude tests generally seek to predict the student's probable strength and weakness in learning a foreign language by measuring performance in an artificial language.

D. Forms of Test

There are two kinds of form of test: objective and subjective test. The distinction between both tests is concern on method of scoring, and nothing else.²⁶ The following explanation will clarify enough about them.

1. Objective Test

Sudijono claimed that objective test is one type of test that is created using items tests, then what the entire test taker has to do is just answering the

²⁴ David P Harris, *Testing English as a second Language*, (New York: Mc Craw-ill, Inc, 1959), Page 3

25 J. B. Heaton, *Writing English Language Test*, (New York: Longman Group, 1988). Page

²⁶ Ibid, Page 22

question by choosing one among several probably answers available in each items or writing sentences or particular symbols in place provided in each item test.²⁷

In line of that, objectives test as cited from Lado is:

“Objectives test are those that are scored rather than mechanically without need to evaluate complex performance on scale”²⁸

a. Types of an objectives test

Sudijono also added that there are five types of objectives test including: true or false test, matching test, completion test, fill in test and also multiple choices. However in this thesis only will clarify the last one.

Multiple choices as stated by Sudijono are a test which is created likely incomplete sentences and the testee should complete the sentence in order to answer the question.²⁹ Before going to design multiple choice test, the test maker or in this case is teacher should know primarily several terms used in multiple choices. First is *stem* which refers to initial part of each multiple choice items. Second is option/responses/alternatives, refers to the options which are available for student to select their answer. One

²⁷ Prof. Drs. Anas Sudijono, *Pengantar Evaluasi Pendidikan*, (Jakarta: PT. Rayagrafindo Pustaka, 1990), Page 106

²⁸ Robert Lado, *Language Testing*, (London: Longman Group, 1961), Page 28

²⁹ Prof. Drs. Anas Sudijono, *Pengantar evaluasi Pendidikan*, (Jakarta: PT. RajaGrafindo Pustaka, 1996), Page 106

2. Subjective Test

As quoted from Lado, subjective test is:

“Tests that require an opinion and a judgment on the part of the examiner”³⁵.

In the other word, Nurgiyantoro have said that subjective test is a test that require student to answer in essay using their word.³⁶

a. Scoring an essay test

Scoring an essay test generally based on the weight of each item test, the level of difficulty, and the amount of the element contained by the answer which is considered as the rightes answer.

For example, there are 5 items test in essay test. The tester had determined that all items have the same level of difficulty, and the elements in each item had made in the same amount. Based on that, tester decided that testee who could answer with the rightest answer or which the answer provides the entire element that required by the tester within the item test, will get 10 marks. When the testee answer almost perfectly or the answer provide mostly the element that required by the test taker, will get 9 mark, and so on.³⁷

³⁵ Robert Lado, *Language Testing*, (London: Longman Group, 1961), Page 28

³⁶ Burhan Nurgiyantoro, *Penilaian dalam Pengajaran Bahasa dan Sastra*, (Yogyakarta, BPFE Yogyakarta, 2001), Page 71

³⁷ Prof. Drs. Anas Sudijono, *Pengantar Evaluasi Pendidikan*, (Jakarta, PT. Raja Grafindo persada, 1996), Page 301

1). The Benefit and the weakness of subjective test

The characteristics of subjective can be seen from its benefits and weakness as follow

a). The Benefit

- Subjective test can create easily and fast
- Avoid students being speculative in answering the items test
- The test taker is able to know how far students understand the material
- Motivate student to organize their thoughts

b). The weakness

- Less able to represent all materials
- It is difficult to score the subjective test. It because the answer of each item might be varieties and wide. Thus, it needs a lot of time, and thoughts to score it.
- Enable test taker to score subjectively
- Validity and reliability of subjective test is poor.

E. Characteristic of A Good Test

All good tests possess three qualities: validity, reliability, and practicality.³⁸ In the other word we say, any test that we use must be appropriate in terms of our objectives (validity), dependable in the evidence it provides (reliability), and applicable to our particular situation (practicality). Without any one of them, a test would be a poor investment in time and money.

Whether the teacher is constructing his own test or selecting a standard instrument for use in his class or school, he should certainly understand what these concepts mean and how to apply them.

1. Validity

Validity of a test is the extent to which it measures what it is supposed to measure.³⁹ In different word but still have same meaning; Gronlund said that validity refers to the appropriateness of interpretations made from test scores and other evaluation result, with regard to a particular use.⁴⁰

Furthermore, a test has validity evidence if we can demonstrate that it measures what it says it measures. For example, if the test is supposed to be a test of tenth-grade English language ability, it should measure tenth-grade

³⁸ David P Harris, *Testing English as a second Language*, (New York: Mc Craw-ill, Inc, 1959), Page 159

J.B Heaton, *Writing English Language Test*, (New York: Longman Group, 1988). Page 172

⁴⁰ Norman E. Gronlund, *Measurement and Evaluation in Teaching*, (New York, collier Macmillan Publisher, 1985), page 55

measure.⁴⁴ Furthermore, the content validity evidence for a test is established by examination. Test questions are inspected to see whether they correspond to what the user decides should be covered by the test.⁴⁵ For instance, when the test is supposed to test some particular objectives of arithmetic material that are taught in third-grade, but in the real it tests out of those objectives or in the other word, it tests other material objectives of arithmetic which are not taught in that class, then the test is called has poor content validity.

However, a test can sometimes look valid but measure something entirely different than what is intended. Content validity is, therefore, more a minimum requirement for a useful test than it is a guarantee of a good test.

In summary, content validity evidence answers the question “Does the test measures the instructional objectives?” In the other word, a test with good content validity evidence matches or fits the instructional objectives.⁴⁶

⁴⁴ Ibid, Page 58.

⁴⁵ Tom Kubiszyn and Gary Borich, *Educational Testing and Measurement* (Singapore, John Wiley & Sons, INC, 2003), Page 300

⁴⁶ Ibid.

pictures that are representative of the objective materials, and also the quality of the color.

Furthermore, a test that has poor face validity is not providing evidence to judge that it is lack of validity, because face validity is not a scientific notion.⁵² Nevertheless, face validity is very important to be considered of writing tests. A test which does not have face validity may not be accepted by the testee, teachers, education authorities or employers.

d. Construct validity

A test has construct validity evidence if its relationship to other information corresponds well with some theory.⁵³ A theory is simply a logical explanation or rationale that can account for the interrelationship among a set of variables.

Construct validity evidence is more specific and immediately practical uses than the others, we may wish to interpret test scores in terms of their psychological meaning.⁵⁴

⁵² Ibid. Page 33

⁵³ Tom Kubiszyn and Gary Borich, *Educational Testing and Measurement* (Singapore, John Wiley & Sons, INC, 2003), Page 302

⁵⁴ Norman E. Gronlund, *Measurement and Evaluation in Teaching*, (New York, collier Macmillan Publisher, 1985), page 72

$$FV = \frac{R}{N}$$

R represents the number of correct answer and N is the number of students taking a test.

Amount the number of index difficulty between 00 to 1,0. These numbers express the difficulty level of an item test. Some experts are difference in giving the amount of number to express the level of difficulty of an item. Oller in Burhan stated that an item test which has index of difficulty between 0,15 up to 0, 85 is adequate, out of that numbers is too easy and too difficult. Thus it is needed to be revised or changed.⁶² However, Arikunto given a common ukuran about the number of index of difficulty as follow⁶³:

- The item test with index of difficulty 00 up to 0,30 means it is difficult
- The item test with index of difficulty 0, 30 up to 0, 70 means it is good
- The item test with index of difficulty 0,70 up to 1,00 means it is easy

⁶² Burhanudin Nurgiyantoro, *Penilaian dalam Pengajaran Bahasa dan Sastra*, (Yogyakarta, BPFE Yogyakarta, 2001)Page 138

212 ⁶³ Suharsimi Arikunto, *Dasar-dasar Evaluasi Pendidikan*, (Jakarta, Bumi Aksara, 1993), Page

2. Index of Discrimination

Index of discrimination of an item is the ability of item test to differentiate the up level student from down level.⁶⁴

The number which shows amount of index of discrimination is called index of discrimination. Like index of difficulty, the number of index of discrimination is between 00 up to 1, 00. However, unlike index of difficulty, index of discrimination has (-) negative sign. It is used when test discriminate in entirely in wrong way showing the quality or ability of testee.⁶⁵ It shows when none of upper level students got a correct answer and the lower level student can answer correctly.

In addition, Heaton stated several ways in analyzing index of discrimination of test item:

1. Arrange the script in rank order of total score and divide into two groups of equal size (the top half and the bottom half). If there is an odd number of script, dispense with one script chosen at random
2. Count the number of those candidate in the upper group answering the first item correctly, then count the number of those candidate in lower level group candidates answering the item correctly, and so on.
3. Subtract the number of correct answers in the upper level group from the number of correct answer in lower level group. Find the difference in the

⁶⁴ Ibid, Page 213

⁶⁵ Suharsimi Arikunto, *Dasar-dasar Evaluasi Pendidikan*, (Jakarta, Bumi Aksara, 1993), Page

4. Divide this difference by the total number of candidates in one group:

(D = Discrimination index, n = Number of candidates in one group, U = Upper level half and L = Lower level half. The index D is the difference between the proportion passing the item in U and L)

Furthermore, Arikunto claimed that item test is called good when it has index of discrimination in range of 0,40 up to 0,70. He also added the classification of the number range of index of discrimination of an item test as follow⁶⁶;

D : 0,70 up to 1,00 means it is excellent

Analysis distractor means analyzing on the distributing of test taker in determining the option of answering the question in the multiple choice test.⁶⁷

⁶⁶ Suharsimi Arikunto, *Dasar-dasar Evaluasi Pendidikan*, (Jakarta, Bumi Aksara, 1993), Page 221

⁶⁷ Ibid, Page 225

who choose the option a, b, c, d, or e, or test taker who do choose none of the option. In this case, called omit (O).

Furthermore, information conducted from analyzing distractor planned to know whether the distractor play good part in the option or not. Distractors must attract more students in lower group. Therefore, if the distractors chosen by more able students, it means they are poor. Besides, distractor which is not chosen by all students shows that they cannot perform well, thus all alternatives must be selected by the test taker. Moreover, when there is the same amount of the voter from both better and poorer students who chose those distractors, means they are still desirable adequate to be used for future test. But when there is divergence of those matters above, the distractors are suggested to be revised and cannot be applied for others test.⁶⁸

In addition, distractors called to be good distractor if it is chosen by at least 5% of students. It also called effective if the omit is chosen by not more than 10%.⁶⁹

⁶⁸ Burhan Nurgiyantoro, *Penilaian dalam Pengajaran Bahasa dan Sastra*, (Yogyakarta, BPFE Yogyakarta, 2001). Page 144.

226 ⁶⁹ Suharsimi Arikunto, *Dasar-dasar Evaluasi Pendidikan*, (Jakarta, Bumi Aksara, 1993), Page

CHAPTER III

RESEARCH METHODOLOGY

In this chapter will employ the research methodology. Means here the writer will formulate the research design that used by the writer in the way of analyzing the study, they are: Research design, Data and Source of Data, Data Collection, Instrument of Data Collection and Data Analysis and Schedule of Research.

A. Research Design

Mardalis categorized four types of research method which are often used; they are Historical research, explorative research, descriptive research and explanatory research.¹ Descriptive research is used to obtain information concerning the current status of the phenomena to describe "what exists" with respect to variables or conditions in a situation. The methods involved range from the survey which describes the status quo, the correlation study which investigates the relationship between variables, to developmental studies which seek to determine changes over time.²

Based on the statement above, the study entitled “Content Validity and Item Analysis on UAS Test for Tenth (X) Grade of SMAN 3 Sidoarjo”, used descriptive research as the way or technique to do the research since it will describe the validity of the content and will also describe the index of difficulty and index of discrimination on UAS test for Tenth (X) Grade of

¹ Drs. Mardalis, *Metode Penelitian*, (Jakarta: Bumi Aksara, 1995). Page 25

² *Ibid.* Page 25

2. The student's answer sheets of The Second Semester Final English Test (2009/2010) for the Tenth grade student of SMAN 3 Sidoarjo
3. The student's score of the Second Semester Final English Test (2009-2010)
4. The standard and basic competencies of 2006 English curriculum for tenth grade of senior high school.

D. Data Collection Techniques

The data in this study will be collected by study documentation. Documentation is a method to get anything on the form of notes, transcripts, magazines, books, etc,⁴ and then data by documentation will be collected through these following steps:

1. Finding the test items, key answer, the student's answer sheets, the teaching materials and student's score of the final test.
2. Finding the Standard and Basic Competencies of the 2006 English curriculum for the tenth grade of senior high school

E. Instrument of Data Collection

According to Mardalis instrument by means of researching is the implement measured. That is with instrument this research could be gathered by the data as the implement to state the mulberry or the percentage as well as more the shortage in the form of quantitative or qualitative⁵

⁴ *Ibid.* Page 234

⁵ Drs. Mardalis, *Metode Penelitian*, (Jakarta: Bumi Aksara, 1995). Page 60

This study used study documentation to measure the validity of collected data and also interview the teachers to support the data.

F. Data Analysis

The result of collected data then will be analyses by using descriptive, means that data will be described as the way it is.

1. Analyzing The Content Validity

In analyzing the content validity, the writer will collect it through the following steps:

- a. Making a list of the standard competencies, basic competencies, indicators, and learning experience for the tenth grade students of senior high school and the indicators of basic competencies given by SMAN 3 for tenth grade student.
- b. Placing each of the test items in the appropriate place with the standard competencies and basic competencies to identify whether or not the standard competencies and basic competencies covered by the final test.
- c. Counting the percentage of the test items of every language aspects.
- d. Concluding the result of analysis.

In order to make these procedures clearer, the writer presents the illustration of the procedures as follows:

Table 1. The example of analyzing Content Validity

STANDARD COMPETENCE	BASIC COMPETENCE	INDICATORS	Learning Experience	ITEM TEST	Σ	%
MENDENGARKAN	Merespon makna dalam percakapan transaksional (to get thing done) dan interpersonal (bersosialisasi) resmi dan tak resmi secara akurat, lancar dan berterima yang menggunakan ragam bahasa lisan sederhana dalam berbagai konteks kehidupan sehari-hari dan melibatkan tindak tutur: berterima kasih, memuji, dan mengucapkan selamat	<ul style="list-style-type: none"> ❖ Mendidentifikasi kata yang didengar, makna kata, hubungan antar pembicara,. ❖ Mengidentifikasi makna tindak tutur, berterimakasih, memuji, mengucapkan salam dan konek situasi. ❖ Merespon tindak tutur, berterimakasih, memuji, ucapan, selamat 	<ul style="list-style-type: none"> • Mendengarkan berbagai tindak tutur yang didengar melalui tape atau teman • Mendiskusikan berbagai tindak tutur yang didengar melalui tape atau teman 	2,5,10,7, 55, 4, 6, 8, 11, 20, 21, 24, 28, 30, 34, 35	5 11	9% 20%

	D	1	4	Good
	O	0	0	NF

SMAN 3 Sidoarjo. There are 6 columns in that table. First column contains standard competence, second column contains basic competencies, the third column contains several indicators that represent the basic competence of the lesson, and the fourth column contains the learning experience or materials that are taught. While the fourth, fifth and the last column are contain the items test that appropriate the basic competence, the sum of the items test that appropriate the basic competence and the percentage of total numbers of particular items that represent the related basic competence.

In addition, according to Nurgiyantoro, the test has the content validity if it covers all the contents as stated in the curriculum. Based on the result of analyzing content validity in table 3 appendix I, the percentage of every aspect of the learning content is conclude as follows:

1. There are 3 or 5% items for speaking which focus on direct and indirect speech explanation
2. There are 3 or 5% items for speaking which focus on passive voice explanation
3. There are 35 or 63% items for reading skill out of 55, which 2 items or 3 % focus on “ Membaca dan memahami pengumuman/surat” lesson, one or 1% focus on “ Memahami makna teks yang dibaca”, 18 items or 32% focus on “Memahami isi teks yang dibaca”, 3 items or 5% focus on “Memperhatikan

dan menemukan cirri-ciri kebahasaan teks” lesson, 4 items or 7% focus on “melengkapi teks dengan kata kerja yang tepat” learning experience, and 7 items or 12% focus on “Menentukan orientasi cerita dengan metode yang berbeda” Learning experience

4. There are also 14 items or 25% out of 55 items test that did not cover the available materials.

Based on the result above, we can conclude that the Semester II English Final Test for Tenth Grade of SMAN 3 Sidoarjo is good since 72% items test represented all materials. Here the agreement of the test is more than 50%. According to Bloom, if the agreement of the test is 50% or more, it can be conclude that the test has high content validity. ¹

Nevertheless, there are still 14 items test or 25% out of all items test did not cover the materials, they items number 8, 12, 17, 18, 20, 21, 22, 25, 45, 48, 49, 52, 53, and 54. Although their content is suitable with the indicators of standard and basic competencies but they were not taught in the class. Those items cannot be tested because students absolutely cannot answer that items test.

¹ Benjamin S. Bloom, at all, *Evaluation to Improve Learning*, (USA, 1981)

B. Analyzing Index of difficulty

To conduct index difficulty value, divide the number of students got the correct answer by the number of students taking test.

In analysing the index of difficulty, first of all, the writer arranged the students' score from the highest score to the lowest one. Then the writer found the top score and the bottom score and divided it into two groups, upper level group and lower level group of equal size.

In addition, the writer treated differently in dividing class into two groups in order to determine upper and lower level in each class in the same amount moreover in X1 and X3 classes. This is because there were odd numbers of script in those classes.

In X1 class there are 3 students belong to upper group who got the same score, 24. It made difficult in dividing the class into two groups in equal size. It should have been 18 students for upper group and 18 students for lower group since the class has 36 students in all, but it became 21 students for upper group and 15 students for lower group. Hence, the writer deleted two number of script in random in order to make it balance as Heaton stated in his book:

"If there is an odd number of a script, dispense with one script chosen in random"²

² J. B Heaton, *Writing English Language Test*, (New York: Longman Group, 1988). Page 178

Although Heaton suggested dispensing with one script if there is an odd number, but in this case the writer deleted two numbers at all. It is because if only one script was deleted, it still does not make sense since the odd number existed in the same order.

Thus, the number of students taking test in X1 that used to compute index of difficulty and index of discrimination later are 34 students. While in X3 class, the writer dispensed with one script since this class has odd number of script. Like stated above, X3 class has 35 number of student taking test. Hence, in this case, the number of test takers that used to calculating index off difficulty as well as index of discrimination is 34 students.

After determining the upper and lower group of students, the writer computed the index of difficulty using Heaton formula as stated in Chapter II and III. Furthermore, the writer used a table to make calculation easy and efficient. The table can be seen in appendix 5-7.

There are six columns in the table. First column contains the number of item. Second column contains the score of students in upper group who answer correctly of each item. Third column contains the score of students in lower group who answer correctly of each item. The fourth and fifth column contains the value of index of difficulty Index and the value of index of discrimination. And the last

column contains the comment. The column of comment divided into comment for Index difficulty value and for index of discrimination value.

However, the writer calculated index of difficulty value of items test for tenth grade of SMAN 3 Sidoarjo. The writer computed the item test for three classes since this thesis used three samples such as X1, X2, and X3. The writer started calculating from X1, X2, and the last X3.

1. Analyzing Index of Difficulty on UAS Test for X1 Class

The result of analysis shown in table 7 appendix 5, reported that there are 27 out of 55 items or 13. 25% of item test have index of difficulty value between 0.32-0.65, they are items number 3, 5, 6, 9, 10, 11, 14, 15, 16, 17, 19, 21, 24, 25 26, 28, 30, 34, 37, 38, 39, 40, 43, 45, 46, 48, and 52. According to Arikunto, the items that show index of difficulty value between 0.30-0.70 are good³. It means that those are categorized adequate items and could be safely used in the future tests without being rewritten. Besides, the items number 12, 18, 22, 23, 29, 31, 32, 33, 35, 36, 47, 49, 50, 51, 53, 54, and 55 are difficult items test since they show difficulty value between 0.03-0.29, as Arikunto said that the items which have index of difficulty value between 0-00-0.30 are difficult item test. Thus, they are needed to be revised because they might desperate students to study more. While, the rest items are easy, such item

³ Suharsimi Arikunto, *Dasar-dasar Evaluasi Pendidikan*, (Jakarta, Bumi Aksara, 1993), Page 212

2. Analyzing Index of Difficulty analysis on UAS Test for X2 Class

[illegible]

previous chapter. Furthermore the writer applied Arikunto's classification to interpret index of discrimination value of each item test value, as follow:⁶

D : 00 up to 0,20 means it is poor

D : 0,20 up to 0,40 means it is satisfactory

D : 0,40 up to 0,70 means it is good

D : 0,70 up to 1,00 means it is excellent

However, the writer started analyzing from X1 class, X2, and then X3 class

1. Analyzing Index Discrimination on UAS Test for X1 Class

The result of analyzing index of discrimination of items test used in X1 as stated in the table 7 appendix 5 recorded that the big number of items tests in X1, it is 15 or 4.09 % of all items test, are poor since the most of items test have value discrimination between 0.06-0.24. They are number 1, 2, 15, 18, 20, 22, 23, 25, 27, 29, 30, 32, 35, 37, 40, 43, 47, 49, and 55. As Arikunto said that poor items dealing with index of discrimination must be revised because it cannot separating the good students from the bad.

The second huge numbers of items tests are categorized as satisfactory since they own discrimination value around 0.24-0.29. It is happen to 14 out of 55 items test, they are items number 3, 5, 7, 12, 13, 16,

⁶ Suharsimi Arikunto, *Dasar-dasar Evaluasi Pendidikan*, (Jakarta, Bumi Aksara, 1993), Page 221

The table 7 in appendix 5, compiled from the result of the index discrimination analysis also recorded there are only 12 items that categorized good because they have discrimination value between 0.41-0.53, means these items are work properly to differentiate upper students from lower students and they could be utilized for the other test.

The table 8 in appendix 6 reported that the items test used in this class are satisfactory since the big number of items test; it is 18 or 5, 8 % of all items test, they are number 5, 10, 16, 20, 24, 25, 26, 27, 30, 31, 33, 34, 36, 38, 39, 40, 42, 44, 48, and 51, have discrimination value between 0.22-

The table 8 in appendix 6 reported that the items test used in this class are satisfactory since the big number of items test; it is 18 or 5, 8 % of all items test, they are number 5, 10, 16, 20, 24, 25, 26, 27, 30, 31, 33, 34, 36, 38, 39, 40, 42, 44, 48, and 51, have discrimination value between 0.22-

0.39. These items may to be refreshed in order they can functioned well.
Because satisfactory items, as Heaton stated, are weak to discriminate upper level from lower level.

There are also 14 items are poor since they value between 0.06-0.17, they are items number 1, 2, 4, 7, 12, 13, 18, 19, 23, 29, 41, 45, 49, and 53. These items are must be revised because poor items cannot put the upper and lower student in the place they should be. Besides, there are 7 items number 8, 9, 11, 35, 46, and 52, did not function at all because their value is 0.00. These items are desirable to rewrite because they did not work at all and it will waste the time even money.

Moreover, there are 10 items that shown negative sign (-), means these items also must be revised since it entirely distinguished in wrong way. In the other hand, the table 8 in appendix 6 recorded that there are only 10 items that categorized as good items that can differentiate students in upper group from the lower group. They are items number 6, 14, 28, 54, and 55. Those items can be kept for the future test. But, however, it is still disappoint because the good items test related on index discrimination pointed only a few numbers.

3. Analyzing Index Discrimination on UAS Test for X3 Class

The results of index discrimination analysis due to X3 in the table 9 of appendix 7 recorded disappoint report since 18 or 5, 8 % of all items test

shown negative sign (-). It means that those items cannot be tested because they discriminate in false path especially items number 1, 2, 3, 5, 10, 12, 14, 15, 18, 23, 26, 27, 29, 35, 37, 44, 45, 46, and 49. They are crucial to be revised since they can manipulate the result of the test as expected, in the other word, these items can block the test to rich its purpose.

Moreover, 16 items test are poor. They are items number 9, 11, 16, 17, 19, 20, 25, 28, 32, 34, 38, 39, 47, 50, 53, and 54. These items have discrimination value between 0.05-0.17. As Arikunto stated that the items test which show index discrimination value between 0.00-0.20 are recognized as poor. It requires deeply thought to revise those tests because they cannot play proper work to discriminate upper students from lower students. Besides, there are 13 items test that are satisfactory because their facility value between 0.23-0.35. They are number 6, 7, 8, 13, 30, 31, 33, 36, 40, 41, 48, 51, and 52. Although they are satisfactory, but they still need to be refreshed in order they can work well. Because satisfactory items test are less good to facilitate item test to discriminate students. However, there are only 3 items test that have facility value 0.00. Means they do not function at all and extremely they must be deleted or at least revised. They are items number 4, 22, and 43.

D. Analyzing the Effectiveness of Distracters.

The effectiveness of distracters is aimed to know whether the item test could work properly as expected or not. The result of analyzing distracters can be used either to revise the poor items and reference to design next other items test.

The analysis of destructors is done by comparing the number of students in upper group with students in lower group selected the false options given. In addition, the good distracters will manipulate more students whose belong to lower group than students belong to upper group. Thus, if there are more able students chosen the distracters, it means that the item does not function as expected in it must be revised.

The result from computing the effectiveness shown in the table 10-12 .in appendix 8-10, There are 5 columns in that table. First column contains the items number. Furthermore, English Final Test for Tenth Grade of SMAN 3 Sidoarjo has fifty five (55) items. The second column contains the alternative given in each items. In this case, there 5 options are available for each item, it is A, B, C, D, E. one alternative is the key answer and the rest are distracter. And it is added by 'omit' where in this case it is written by 'O'. Omit is used for both upper and lower students who chose nothing the incorrect alternatives.

In addition, the following explanations are used to describe the result of analyzing items distracters which is already written in that table.

adequate distracters are desirable may be used for the future tests without being revised.

2. Analyzing the Effectiveness of Distractors in X2 Class

Based on the result of analyzing the effectiveness of distractors of items test used for X2 shown in the table 11 of appendix 9, it can be seen that there are 130 distractors that good since they are chosen by more bad students than good students. It is relief because they could be kept and applied for others tests. There are also 127 non functioned distractors because no one from upper and lower student voted those alternatives. Thus, those distractors must be revised as Nurgiyantoro have said that if the distractors

Besides, there are 57 distractors are needed to be changed because they are malfunction. Those items are selected by more good students than bad students. However, good distractors should have attracted more students in lower group than students in upper group.

Furthermore, there are 16 distractors that adequate because the good students who chosen those distractors have the same amount with the lower students. According to Nurgiyantoro, those distractors are acceptable to use for future tests.

In conclusion, we simply say that the distractors of items test used for X2 are good since 133 out of 330 distractors performed efficiently to attract more

bad students than lower students. Hence, those distractors can be used for the next test. Nevertheless, this still disappoint since the non functioned distractors also shown big number, it is 104, or almost half distractors. Thus, the test makers or the teachers have to be more aware to revise those distractors.

3. Analyzing the Effectiveness of Distractors for X3 Class

Like in X2, the result of analyzing the effectiveness used for X3 as stated in the table 12 of appendix 10 recorded the similar report. There are 126 good distracters that can perform efficiently to attract more students in lower level. According to Nurgiyantoro, those items could be safely used for other tests.

There are also 107 items that has no function since they are not selected by both able and less able students. Thus, those distractors are not good to apply for other test and must be replaced.

Moreover, there are 78 distractors that worked in contrary as expected since they attracted the wrong candidates (i.e the better ones). Furthermore, there are 19 distractors that adequate because either upper or lower students have the same number in selecting those options.

2. The index of difficulty of Semester II English Final Test for Tenth Grade of SMAN 3 Sidoarjo used for X1 is acceptable 27 or 13, 25% out of 55 items have facility value between 0.32-0.65 or categorized as good items. Thus, those items could safely be used for future test without being rewritten. While the rest items have to be revised because they are too easy and difficult. Moreover, since almost a half test, it is 24, shown difficulty value between 0.31-0.69, the items test used for X2 of SMAN 3 Sidoarjo are acceptable too. It is because their facility value is recognized as good test. For items test that are easy and difficult, they must be changed because they cannot work properly. Unlike X1 and X2, index of difficulty of UAS test used for X3 is easy since 21 out of 55 items test have difficulty value between 0.74-0.91. Hence, those items are not desirable to apply for other test and needed to be revised as well as the items that are difficult.
3. Index of discrimination of UAS test used for X1 is poor since the biggest number of items, it is 15 out of 55, have discrimination value between 0.06-0.24. Those items badly need to be revised as well as the items that have satisfactory, non function, and even miss function index of discrimination. It is because all of those items cannot perform to distinguish better student from poorer student. While the rest items that have good index discrimination do not need to be revised and can be used for future test. Similarly, the index of discrimination of UAS test used for X2 is satisfactory since 18 out of 55

items test shown index discrimination value between 0.22-0.39. It indicates that those items need to be rewritten as well as those items are recognized poor, non function, and malfunction. While, 10 items that are categorized good are could be kept and applied for others tests. On the other hand, the index of discrimination in X3 are known malfunction since the biggest amount or items test, it is 18, have negative (-) sign. It means that those items cannot perform correctly because they discriminate the wrong candidate (better students). Thus, they badly need to be revised.

4. The result of analyzing the effectiveness of distractors of UAS test for X1, X2, and X3 class recorded the same report, that the items test are good since 166 or 81, 51% out of 330 distractors of items test for X1, 130 or 45% out of all distractors of items test for X2, and 126 or out of 330 distractors of items test for X3 are performed efficiently. It means that those distractors could safely be used for future test, while the rest items that cannot work perfectly have to be revised.

B. SUGGESTION

After recognizing the result of this study related to content validity and items analysis of Semester II English Final Test for Tenth Grade of SMAN 3 Sidoarjo, there are several matters that are seemed to be suggested.

As stated in the first chapter, the suggestion due to the teacher, the students, and the further researcher, as follows:

For the teacher:

1. The teachers should try their own test out to know whether that tests are adequate before it is given to student as well as analyze it after it is tested. Hopefully, their test could perform correctly to measure the progress of their students as the function of good test.
2. Although the result of content validity of Semester II English Final Test for Tenth Grade of SMAN 3 Sidoarjo is good, but teachers still should modify the items to be better items test since most of the items only cover one skill, it is reading, and 14 items were not taught in that class, moreover if the items test will be used once more in the future test.
3. Teachers should revise the items test that are poor, and difficult, in order there will no items that too easy and too difficult. Thus, teachers will know the progress of students correctly.
4. The teachers must aware to set the amount of items test for students. In the other word, we can say that amount of items test given must suitable with the time available for doing that test. If the items test are many, the time settled for work it must be added and vice versa. The amount of items test for Semester II English Final Test for Tenth Grade of SMAN 3 Sidoarjo is not

suitable with the time available since the items are 55 and the time is only 90 minutes. It can influence student to rush the time left and much probably they will answer the items without carefully, and

5. Teachers should revise the items test that are poor, satisfactory, non function, mal function and even adequate, in order there will no items that discriminate in wrong way and could distinguish the better students from poorer students. Besides, teachers also should modify the distractors that are not good (non function, malfunction).

For students: Students should be able to recognize which tests those are good to do and which tests those are not good. Because bad test do not benefit them, in the other word, the bad test cannot measure their progress.

For further researcher: there must be several tests that are needed to be researched, in order to repair and fix all tests given and avoid from the designing and using of test that ignore of the criteria of good test (content validity, reliability, index of difficulty, index of discrimination, and the effectiveness of distractors).

BIBLIOGRAPHY

- Arikunto, Suharsimi. 1996. *Prosedur Penelitian*. Jakarta: PT. Rineka Cipta.
- Arikunto, Suharsimi. 1984. *Dasar-dasar Evaluasi Pendidikan*. Jakarta: Bina Aksara.
- Bloom, Benjamin S. et al.. 1981. *Evaluation to Improve Learning*. USA
- Brown, H. Douglas. 2001. *Teaching by Principles: An Interactive Approach to Language Pedagogy*, Second Edition. San Francisco State University: Longman, Inc
- Cholbi, Laiyinatul. 2006. The Final English Test for the fourth grade students of SDN Pucang III Sidoarjo. Unpublished S-1 Thesis. Surabaya: UNESA (FBS)
- Daryanto, Drs. H. 1999. *Evaluasi Pendidikan*. Jakarta: Rineka Cipta.
- Depdiknas. 2004. *Standard Kompetensi Mata Pelajaran Bahasa Inggris SMP dan Madrasah Tsanawiyah*. Jakarta: Departemen Pendidikan Nasional
- Ebel, Robert L. 1979. *Essential of Educational Measurement*. USA: Prentice-Hall INC New Jersey
- Gary Borich, Tom Kubiszy. 2003. *Educational Testing and Measurement*. Singapore: John Wiley & Sons, INC.
- Gronlund, Norman E. 1985. *Measurement and Evaluation in Teaching*. New York: collier Macmillan Publisher.
- Hadi, Sutrisno. 1989. *Metodologi Research*. Yogyakarta: Andi offset.
- Harris, David P. 1959. *Testing English as a second Language*. New York: Mc Craw-ill, Inc.
- Hughes, Arthur. 2003. *Testing for Language Teacher*. Cambridge: University Press.

- Heaton, J. B. 1988. *Writing English Language Test*. New York: Longman Group.
- Lado, Robert. 1961. *Language Testing*, London: Longman Group.
- Lampiran Peraturan Menteri Pendidikan Tentang Sistem Pendidikan Nasional No. 20 tahun 2003
- Mardalis, Drs. 1995. *Metode Penelitian*. Jakarta: Bumi Aksara.
- Nurdiyantoro, Burhan. 2001. *Penilaian dalam Pengajaran Bahasa dan Sastra*. Yogyakarta: BPFE Yogyakarta.
- Purwanto, Drs. M. Ngilim. 1985. *Prinsip-prinsip dan Teknik evaluasi Pengajaran*. Bandung: Remadja Karya.
- Peraturan Republik Indonesia Tentang Sistem Pendidikan Nasional No. 20 tahun 2003
- Pratikasari, Ria Dhewi. 2006. Semester II English Summative Test for The Eighth Year Student of SMP Negeri 1 Slahung Ponorogo. Unpublished S-1 Thesis. Surabaya: UNESA (FBS)
- Rea-Dickins, Pauline& Kevin Germaine. 2008. *Evaluation*. New York: Oxford University Press.
- Sumartana, Wayan Nurkanca, 1986. *Evaluasi Pendidikan*. Surabaya: Usaha Nasional.
- Thoha, Drs. M. Chabib. 1991. *Teknik Evaluasi Pendidikan*. Jakarta: PT Raja Grafindo Persada.